



AnnotQTL: a new tool to gather functional and comparative information on a genomic region

Frédéric Lecerc, Anthony Bretaudeau, Olivier Sallou, Colette Désert, Yuna Blum, Sandrine Lagarrigue, Olivier Demeure

► To cite this version:

Frédéric Lecerc, Anthony Bretaudeau, Olivier Sallou, Colette Désert, Yuna Blum, et al.. AnnotQTL: a new tool to gather functional and comparative information on a genomic region. *Nucleic Acids Research*, 2011, 39 (Suppl 2), open access. 10.1093/nar/gkr361 . hal-00597398

HAL Id: hal-00597398

<https://hal.science/hal-00597398>

Submitted on 3 Jun 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

AnnotQTL: a new tool to gather functional and comparative information on a genomic region

F. Lecerf^{1,2,*}, A. Bretaudeau³, O. Sallou³, C. Desert^{1,2}, Y. Blum^{1,2,4}, S. Lagarrigue^{1,2} and O. Demeure^{1,2}

¹INRA, UMR598 Génétique Animale, F-35000 Rennes, ²Agrocampus OUEST, UMR598 Génétique Animale, F-35000 Rennes, ³GenOuest Platform, INRIA/Irisa – Campus de Beaulieu, F-35042 Rennes Cedex and ⁴Agrocampus OUEST, Applied Mathematics Department, F-35000, Rennes, France

Received February 9, 2011; Revised April 18, 2011; Accepted April 27, 2011

ABSTRACT

AnnotQTL is a web tool designed to aggregate functional annotations from different prominent web sites by minimizing the redundancy of information. Although thousands of QTL regions have been identified in livestock species, most of them are large and contain many genes. This tool was therefore designed to assist the characterization of genes in a QTL interval region as a step towards selecting the best candidate genes. It localizes the gene to a specific region (using NCBI and Ensembl data) and adds the functional annotations available from other databases (Gene Ontology, Mammalian Phenotype, HGNC and Pubmed). Both human genome and mouse genome can be aligned with the studied region to detect synteny and segment conservation, which is useful for running inter-species comparisons of QTL locations. Finally, custom marker lists can be included in the results display to select the genes that are closest to your most significant markers. We use examples to demonstrate that in just a couple of hours, AnnotQTL is able to identify all the genes located in regions identified by a full genome scan, with some highlighted based on both location and function, thus considerably increasing the chances of finding good candidate genes. AnnotQTL is available at <http://annotqtl.genouest.org>.

INTRODUCTION

The final steps of genetic mapping research programs require close analysis of several QTL regions to select candidate genes for further studies. Despite several websites

(NCBI genome browser, Ensembl Browser, UCSC Genome Browser) or web tools (Biomart, Galaxy) developed to achieve this task, the selection of candidate genes remains a laborious process. The information made available on the more prominent web sites differs slightly in terms of gene prediction and functional annotation, while other websites provide extra information that researchers may want to use (HGNC approved gene symbols, Gene Ontology (GO) Annotation or functional data, conservation of synteny with other species, etc.). It is possible to manually merge and compare this information for one QTL containing few genes, but not for many different QTL regions containing dozens of genes.

Here, we propose a web tool that, for a given region of interest, merges the list of genes available in NCBI and Ensembl, removes redundancy, adds functional annotations from different prominent web sites, and highlights the genes for which functional annotation fits the biological function or diseases of interest. The tool is dedicated to sequenced species of livestock including cattle, pig, chicken and horse as well as dog, i.e. species that have been extensively studied (with over 8000 QTLs detected; see <http://www.animalgenome.org/cgi-bin/QTLdb/index>). Nevertheless, because of the family designs and the low number of animals used in these species, most of the studies use linkage analysis, and the QTL regions identified remain large (containing dozens of genes). Conversely, in human and model species, most analyses now draw heavily on association studies involving large cohorts, thus providing more power and accuracy, and the web tools already available focus on these species through functional annotation of SNPs in association with the trait (1–8). Most of these tools focus on the SNP annotation itself, describing whether the SNP is located in a gene, or even in a coding sequence, and defining if it could have a functional effect. While these web tools are highly efficient in providing a good annotation for specific SNPs, they

*To whom correspondence should be addressed. Tel: +(33) 2 23 48 59 62; Fax: +(33) 2 23 48 54 70; Email: lecerf@agrocampus-ouest.fr

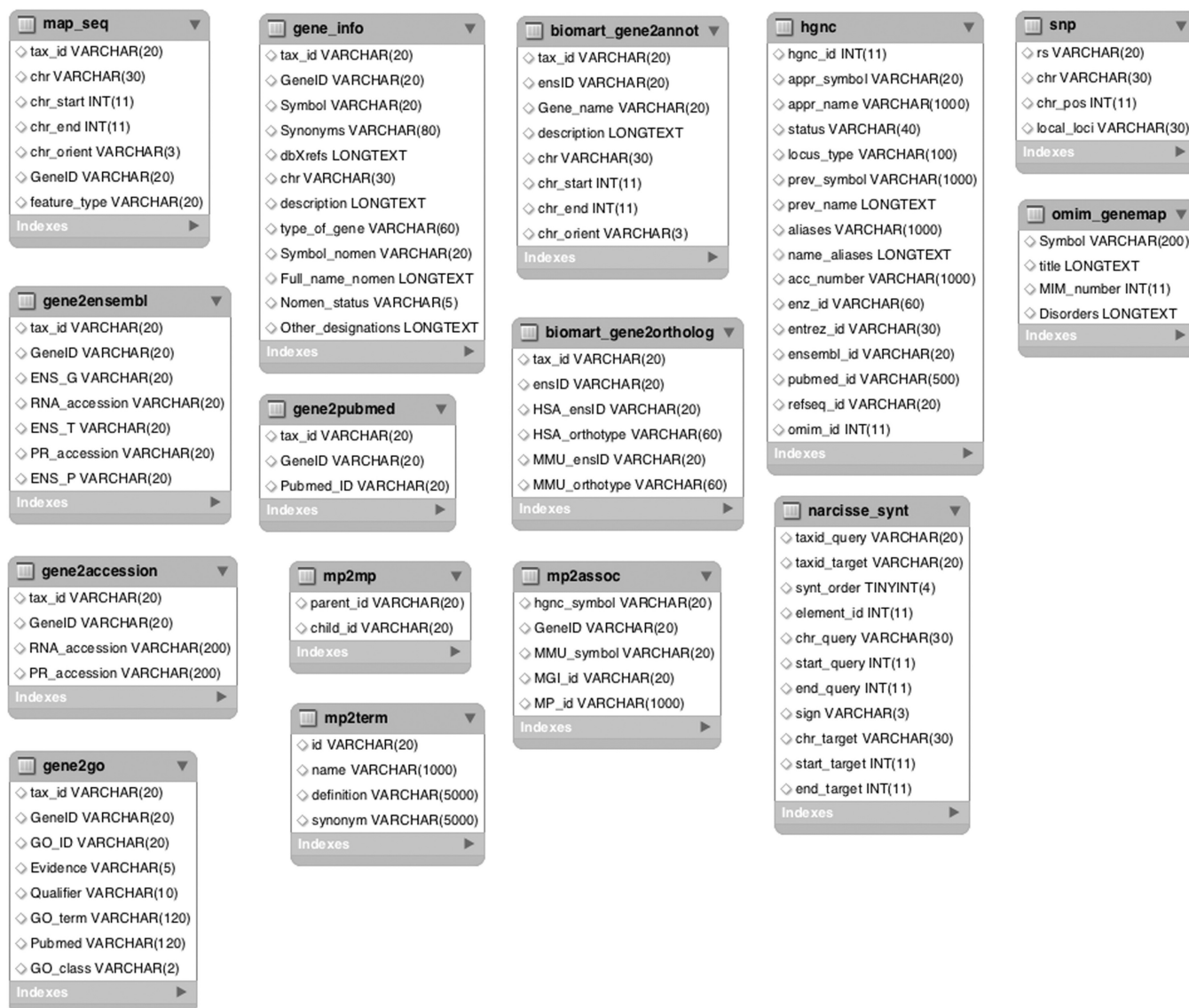


Table 1. Statistics of the QTL/eQTL regions analyzed using AnnotQTL

| | Number of regions | Regions mean size (Mb) | NCBI genes | Ensembl genes | AnnotQTL genes obtained merging NCBI and Ensembl | GO and MP terms screening results | |
|------|-------------------|------------------------|------------|---------------|--|-----------------------------------|-----------------------------------|
| | | | | | | Genes found | Average of genes found per region |
| QTL | 21 | 4.8 | 1734 | 1902 | 2220 | 127 | 5.8 |
| eQTL | 25 | 3.4 | 1198 | 1283 | 1506 | 93 | 3.7 |

working from an initial set of 1734 genes from the NCBI database and 1902 genes from the Ensembl database, AnnotQTL retrieved a non-redundant set of 2220 genes. On this large dataset, we applied the ‘highlight function’ on each region to underline genes whose functional annotation was related to the studied phenotype. Among the 2220 genes located in these 21 QTL regions, 127 were highlighted using the GO term, ‘lipid’ and the MP term ‘adipose’ as keywords, with an average 5.4 genes highlighted per region.

Finally, AnnotQTL can also be exploited to look at eQTL regions. Strategies combining transcriptomics and genotyping data have recently been developed to better characterize QTL regions for traits of interest by identifying co-localized eQTLs and QTLs (16–21). Whatever the context, this strategy identifies a much higher number of eQTL regions than in QTL studies, thus creating a need for tools that can efficiently find positional and functional candidate genes. Here, we focus on 25 chicken eQTL regions affecting 70 genes involved in lipid metabolism (i.e. sharing the GO term GO:0006629 ‘lipid metabolic process’). Average length of these regions is 3.4 Mb. Running AnnotQTL found similar results to those obtained for the QTL regions. All the regions were enriched with genes by comparing NCBI and Ensembl information against information provided by either NCBI or Ensembl only (Table 1): working from an initial set of 1,198 genes from the NCBI database and 1283 genes from the Ensembl database, AnnotQTL retrieved a non-redundant set of 1506 genes. Again, in order to select possible candidate genes, we used the ‘highlight function’ to pinpoint the genes related to the studied phenotype. Among these 1506 genes, and using the same GO term ‘lipid’ and MP term ‘adipose’ as keywords, a total of 93 genes were identified, with an average 3.7 genes highlighted per region.

These examples corresponding to two different contexts (QTL and eQTL analyses) clearly demonstrate how in just a couple of hours, AnnotQTL can accurately analyze the gene content of numerous regions identified by a full genome scan and go on to highlight some of these genes based on both their location and function, whereas in the same time period, a manually run procedure would only have been able to analyze one single region.

CONCLUSION

AnnotQTL is a web tool designed to gather the functional annotation of different prominent web sites while minimizing redundant information. Using all known

information substantially accelerates the gene analysis of QTL regions for livestock species traits and improves the selection of candidate genes.

ACKNOWLEDGEMENTS

The authors thank A.T.T. scientific editing services for proofreading the article and all the beta testers for their help on debugging AnnotQTL.

FUNDING

INRA, Agrocampus Ouest, the Regional Council of Brittany; French Ministry in charge of Agriculture (DGER). Funding for open access charge: INRA.

Conflict of interest statement. None declared.

REFERENCES

- Yue, P., Melamud, E. and Moul, J. (2006) SNPs3D: candidate gene and SNP selection for association studies. *BMC Bioinformatics*, **7**, 166–166.
- Goodswen, S., Gondro, C., Watson-Haigh, N. and Kadarmideen, H. (2010) FunctSNP: an R package to link SNPs to functional knowledge and dbAutoMaker: a suite of Perl scripts to build SNP databases. *BMC Bioinformatics*, **11**, 311.
- Wang, P., Dai, M., Xuan, W., McEachin, R. C., Jackson, A. U., Scott, L. J., Athey, B., Watson, S. J. and Meng, F. (2006) SNP Function Portal: a web database for exploring the function implication of SNP alleles. *Bioinformatics*, **22**, e523–e529.
- Shen, T. H., Carlson, C. S. and Tarczy-Hornoch, P. (2009) SNPit: a federated data integration system for the purpose of functional SNP annotation. *Comp. Methods Prog. Biomedicine*, **95**, 181–189.
- Ryan, M., Diekhans, M., Lien, S., Liu, Y. and Karchin, R. (2009) LS-SNP/PDB: annotated non-synonymous SNPs mapped to Protein Data Bank structures. *Bioinformatics*, **25**, 1431–1432.
- Riva, A. and Kohane, I. S. (2004) A SNP-centric database for the investigation of the human genome. *BMC Bioinformatics*, **5**, 33.
- Reumers, J., Maurer-Stroh, S., Schymkowitz, J. and Rousseau, F. (2006) SNPeff v2.0: a new step in investigating the molecular phenotypic effects of human non-synonymous SNPs. *Bioinformatics*, **22**, 2183–2185.
- Li, S., Ma, L., Li, H., Vang, S., Hu, Y., Bolund, L. and Wang, J. (2007) Snap: an integrated SNP annotation platform. *Nucleic Acids Res.*, **35**, D707–D710.
- Haider, S., Ballester, B., Smedley, D., Zhang, J., Rice, P. and Kasprzyk, A. (2009) BioMart Central Portal—unified access to biological data. *Nucleic Acids Res.*, **37**, W23–W27.
- Courcelle, E., Beausse, Y., Letort, S., Stahl, O., Fremez, R., Ngom-Bru, C., Gouzy, J. and Faraut, T. (2008) Narcisse: a mirror view of conserved syntenies. *Nucleic Acids Res.*, **36**, D485–D490.
- Filangi, O., Beausse, Y., Assi, A., Legrand, L., Larre, J. M., Martin, V., Collin, O., Caron, C., Leroy, H. and Allouche, D. (2008) BioMAJ: a flexible framework for databanks synchronization and processing. *Bioinformatics*, **24**, 1823–1825.

12. Charlier,C., Coppeters,W., Farnir,F., Grobet,L., Leroy,P.L., Michaux,C., Mni,M., Schwes,A., Vanmanshoven,P., Hanset,R. *et al.* (1995) The mh gene causing double-muscling in cattle maps to bovine Chromosome 2. *Mamm. Genome*, **6**, 788–792.
13. Grobet,L., Martin,L.J., Poncelet,D., Pirottin,D., Brouwers,B., Riquet,J., Schoeberlein,A., Dunner,S., Menissier,F., Massabanda,J. *et al.* (1997) A deletion in the bovine myostatin gene causes the double-muscléd phenotype in cattle. *Nat. Genet.*, **17**, 71–74.
14. Lagarrigue,S., Pitel,F., Carre,W., Abasht,B., Le Roy,P., Neau,A., Amigues,Y., Sourdioux,M., Simon,J., Cogburn,L. *et al.* (2006) Mapping quantitative trait loci affecting fatness and breast muscle weight in meat-type chicken lines divergently selected on abdominal fatness. *Genet. Sel. Evol.*, **38**, 85–97.
15. Le Mignon,G., Pitel,F., Gilbert,H., Le Bihan-Duval,E., Vignoles,F., Demeure,O., Lagarrigue,S., Simon,J., Cogburn,L.A., Aggrey,S.E. *et al.* (2009) A comprehensive analysis of QTL for abdominal fat and breast muscle weights on chicken chromosome 5 using a multivariate approach. *Anim. Genet.*, **40**, 157–164.
16. Blum,Y., Le Mignon,G., Lagarrigue,S. and Causeur,D. (2010) A factor model to analyze heterogeneity in gene expression. *BMC Bioinformatics*, **11**, 368.
17. Le Mignon,G., Desert,C., Pitel,F., Leroux,S., Demeure,O., Guernec,G., Abasht,B., Douaire,M., Le Roy,P. and Lagarrigue,S. (2009) Using transcriptome profiling to characterize QTL regions on chicken chromosome 5. *BMC Genomics*, **10**, 575.
18. Ponsuksili,S., Jonas,E., Murani,E., Phatsara,C., Srikanchai,T., Walz,C., Schwerin,M., Schellander,K. and Wimmers,K. (2008) Trait correlated expression combined with expression QTL analysis reveals biological pathways and candidate genes affecting water holding capacity of muscle. *BMC Genomics*, **9**, 367.
19. Babak,T., Garrett-Engel,P., Armour,C.D., Raymond,C.K., Keller,M.P., Chen,R., Rohl,C.A., Johnson,J.M., Attie,A.D., Fraser,H.B. *et al.* (2010) Genetic validation of whole-transcriptome sequencing for mapping expression affected by cis-regulatory variation. *BMC Genomics*, **11**, 473.
20. Drost,D.R., Benedict,C.I., Berg,A., Novaes,E., Novaes,C.R., Yu,Q., Dervinis,C., Maia,J.M., Yap,J., Miles,B. *et al.* (2010) Diversification in the genetic architecture of gene expression and transcriptional networks in organ differentiation of *Populus*. *Proc. Natl Acad. Sci. USA*, **107**, 8492–8497.
21. van Nas,A., Ingram-Drake,L., Sinsheimer,J.S., Wang,S.S., Schadt,E.E., Drake,T. and Lusis,A.J. (2010) Expression quantitative trait loci: replication, tissue- and sex-specificity in mice. *Genetics*, **185**, 1059–1068.